1

**ENTRAIN: integrating trajectory inference and gene regulatory networks with spatial data to co-localize the receptor-ligand interactions that specify cell fate**

Wunna Kyaw[1,2], Ryan C. Chai[2,3], Weng Hua Khoo[2,3], Leonard D. Goldstein[4,5], Peter I. Croucher[2,3], John M. Murray[6], Tri Giang Phan[1,2]

**Affiliations**: [1]Precision Immunology Program, Garvan Institute of Medical Research, Darlinghurst, NSW, Australia; [2]St Vincent's Healthcare Clinical Campus, Faculty of Medicine and Health, UNSW Sydney, Darlinghurst, NSW, Australia; [3]Skeletal Diseases Program, Garvan Institute of Medical Research, Darlinghurst, NSW, Australia; [4]Kinghorn Centre for Clinical Genomics, Garvan Institute of Medical Research, Darlinghurst, NSW, Australia; [5]St George and Sutherland Clinical School, UNSW Medicine and Health, UNSW Sydney, Darlinghurst, NSW, Australia; [6]School of Mathematics and Statistics, Faculty of Science, UNSW Sydney.

**Correspondence**: Dr Tri Phan, Garvan Institute of Medical Research, 384 Victoria St, Darlinghurst, NSW 2010, Australia; Email: t.phan@garvan.org.au

**Keywords**: trajectory inference; ligand-receptor pair; cell-cell communication; cellular niche; cell fate; single cell transcriptomics; spatial transcriptomics; random forest

27 **Abstract**

28 Cell fate is commonly studied by profiling the gene expression of single cells to infer

29 developmental trajectories based on expression similarity, RNA velocity, or statistical

30 mechanical properties. However, current approaches do not recover

31 microenvironmental signals from the cellular niche that drive a differentiation

32 trajectory. We resolve this with environment-aware trajectory inference (ENTRAIN),

33 a computational method that integrates trajectory inference methods with ligand-

34 receptor pair gene regulatory networks to identify extracellular signals and evaluate

35 their relative contribution towards a differentiation trajectory. The output from

36 ENTRAIN can be superimposed on spatial data to co-localize cells and molecules in

37 space and time to map cell fate potentials to cell-cell interactions. We validate and

38 benchmark our approach on single-cell bone marrow and spatially resolved

39 embryonic neurogenesis datasets to identify known and novel environmental drivers

40 of cellular differentiation. ENTRAIN is available as a public package at

41 https://github.com/theimagelab/entrain and can be used on both single-cell and

42 spatially resolved datasets.

43

44 **Main text**

45 In multicellular organisms, cells in different organs and tissues adopt different states

46 of cellular differentiation to allow them to perform specialized tasks. The precise

47 coordination of cellular differentiation and function requires not only the existence of

48 multiple distinct cellular fates but also the ability of the cells to communicate and

49 regulate each other to maintain homeostasis and avoid disease[1]. The development

50 of single-cell technologies such as single-cell RNA sequencing (scRNA-seq) has

51 revolutionized our ability to deconvolute the myriad of heterogenous cellular

52 transcriptional states that comprise multicellular life, even in seemingly homogenous

53 cell lineages such as natural killer (NK) cells[2]. Interestingly, scRNA-seq has

54 suggested that cells exist in a continuum of transcriptional states, whereas the

55 traditional assignment of cell identity by the expression of cell lineage markers, such

56 as by flow cytometry, have viewed cell fates as discrete, non-overlapping entities[3].

57 Thus, the cell state is the transcriptional output of the gene regulatory networks and

58 may represent transient intermediate steps in the differentiation of the cell towards its

59 developmental destination, or cell fate[4, 5]. Accordingly, it may also be possible to

60    predict the future cell fate from the current cell state and the dynamic expression of

61    critical master regulator genes.

62

63    Trajectory inference computes the pattern of change in gene expression for cells in a

64    given dataset and arranges them in pseudo-chronological order along a

65    developmental pathway (pseudotime) based on the similarity between their changing

66    gene expression profiles[6, 7]. There are currently more than 70 published trajectory

67    inference methods, with many more in development[6]. This reflects both the

68    popularity of pseudotime for lineage tracing and also the limitations of the technique,

69    which are dependent on the underlying assumptions, many of which are project and

70    cell-type specific[8]. RNA velocity is an alternative approach that uses the relative

71    abundance of unspliced to spliced mRNA transcripts to predict future cell states,

72    instead of inferring them from global similarity in the transcriptomic profiles between

73    cells[9, 10]. However, the modelling of RNA kinetics also makes several assumptions,

74    such as a common rate of splicing across different genes and the sampling of

75    multiple intermediate cell states in addition to the mature steady-state[11]. The RNA

76    velocity analysis of peripheral blood mononuclear cells (PBMCs), which contain

77    mature blood cells without the immature bone marrow precursor cells, is a good

78    example of the potential for this approach to generate spurious cell lineage

79    relationships[11, 12]. Thus, there are fundamental limits to the fidelity of dynamic

80    inferences that can be made from single cell snapshots[13]. The cross-validation of cell

81    state transitions and lineage relationships by additional orthogonal methods has

82    therefore been strongly recommended[11, 12].

83

84    The development of tools for ligand-receptor (LR) network analysis of single cell data

85    has made it possible to decipher the cell-cell communications that may also drive cell

86    state transitions and determine cell fate[1]. First used to infer cellular interactions at

87    the feto-maternal interface in the human placenta[14], LR analysis has become

88    increasingly popular with its ability to infer interactions between cells in a given

89    dataset, even in the absence of spatial information[15]. Broadly, tools for LR analysis

90    can be generalized into two categories: 1. 'LR-only' tools that rely solely on ligand-

91    receptor gene expression, and 2. 'LR + Intracellular' tools that incorporate

92    intracellular regulons. 'LR-only' tools, such as CellPhoneDB[16, 17], predict cell-cell

93    interactions by considering the expression of ligand and receptor genes as a proxy

94   for secreted and membrane protein abundance. Tools from the 'LR + Intracellular'
95   category are motivated by the possibility that a scarcely expressed LR pair may also
96   unexpectedly regulate a considerable array of downstream genes, which would be
97   overlooked by 'LR-only' tools that only consider gene expression levels. To this end,
98   these tools exploit the large body of biological prior knowledge about gene regulatory
99   networks and intracellular signalling pathways to prioritize LR pairs based on their
100  downstream influence on gene regulation.  As a result, methods belonging to the 'LR
101  + Intracellular' category achieve markedly different results from methods in the 'LR-
102  only' category. Thus, LR analysis has potential to complement trajectory inference
103  and RNA velocity by providing corroborating evidence for gene regulatory
104  programmes responsible for cell state transitions. However, only two tools belong to
105  the second category, NicheNet[18] and CellCall [19], and no tools to date incorporate
106  trajectory or velocity information with LR analysis.

107

108  The introduction of spatially resolved transcriptomics has demonstrated the
109  important role of physical location within a tissue. Specifically, different stages of
110  differentiation within a population often correlate with microanatomical location in the
111  tissue[20]. Similarly, LR interactions are limited by surface contact between interacting
112  cells, or through diffusivity for secreted ligands[21]. This suggests that the spatial
113  information of a cell, which is typically lost in traditional scRNA-seq workflows, can
114  improve the evaluation of LR pairs that influence the differentiation trajectories of a
115  cell. Therefore, there is a need for computational methods that incorporate spatially
116  resolved data to better understand the environmental drivers of differentiating
117  populations.

118

119  Here, we have integrated the information provided by trajectory inference and RNA
120  velocity with LR analysis to develop ENTRAIN, an environment-aware trajectory
121  inference computational tool that can be used to predict the extracellular drivers of
122  cell state transitions. ENTRAIN consists of three modules, ENTRAIN-Pseudotime,
123  ENTRAIN-Velocity, and ENTRAIN-Spatial, which can be applied on the outputs of
124  pseudotime-based methods, RNA velocity or paired single-cell and spatially resolved
125  data, respectively. In turn, ENTRAIN can be applied to a wide range of datasets
126  containing differentiating cells as well as the cell's interacting microenvironment,

127 including spatial datasets. The ENTRAIN package is available to download at

128 https://github.com/theimagelab/entrain.

129

130 **METHODS**

131 **Materials and Methods**

132 **Assumptions and Overview**

133 ENTRAIN operates based on certain assumptions about the biological system of

134 interest:

135

136 1) Environmental control over a differentiating cell population, if present, is

137 facilitated through LR interactions.

138 2) The environmental influence on differentiation is operating on a time scale

139 resolvable by either pseudotime-based or RNA velocity methods.

140 3) The environmental regulation occurs via known regulatory pathways that are

141 documented in gene regulatory network databases, and that the degree of

142 regulation in this database can be quantified as the edge weight ($w$) between a

143 given ligand ($l$) and a given gene $g \in G$, where $G$ denotes the set of all genes in

144 the genome.

145

146 The fundamental operating principle of ENTRAIN is that, if a specific ligand $l$ is

147 influencing the expression of a specific gene $g$ in a differentiating population, this

148 influence can be observed as a meaningful contribution of the ligand-gene regulatory

149 network towards predicting the observed changes in the expression of $g$. In other

150 words, if the edge weight $w$ between $l$ and $g$, which represents the strength of the

151 regulatory interaction, positively correlates with the observed gene expression

152 changes in the trajectory (or velocity), then this suggests that the ligand is actively

153 driving the observed differentiation for that gene in the observed dataset.

154

155 First, we construct differentiation trajectories either by using manifold-based

156 trajectory inference tools[7] or RNA velocity estimation with scVelo[10]. We then identify

157 trajectory informative ('TRAINing') genes that either correlate their expression with

158 pseudotime (for manifold-based trajectories) or exhibit high velocity likelihoods (for

159 RNA velocity-based methods). In parallel, we identify LR pairs using NicheNet[18] and

160  extract regulatory interactions between identified LR pairs and downstream target

161  genes in the regulon. We then fit a random forest regression model using TRAINing

162  gene covariances (for pseudotime) or velocity probabilities (for scVelo) as the

163  'response' variable and NicheNet predicted regulatory interactions as the

164  'explanatory' variable. This model estimates the proportion of trajectory dynamics (as

165  measured by pseudotime covariance or velocity likelihood) that can be predicted by

166  the regulatory interactions downstream of a LR pair. Ligands are scored based on

167  their contributions to the model.

168

169  **Trajectory construction with Monocle**

170  Consider cells as $n$ vectors in $\mathbb{R}^{|G|}$, where $|G|$ is the number of genes measured by

171  the scRNA-seq experiment and $n$ is the number of cells. Typically, a differentiation

172  process will take the form of an ordered sequence of cells in this high dimensional

173  space, beginning at a root cell (or node), traversing along a series of intermediate

174  cells with progressive changes in gene expression before ending at a terminal cell. In

175  this ordered sequence, called pseudotime, cells that are highly similar in gene

176  expression space will be adjacent in pseudotime. Assuming sufficient sampling of

177  intermediate cell stages, this approach successfully identifies differentiation

178  trajectories but cannot determine whether a trajectory is driven by its environment or

179  is under cell-intrinsic control, motivating the use of ENTRAIN to identify

180  environmental influences. ENTRAIN implements pseudotime analysis by using the

181  Monocle3[22] workflow, which applies the SimplePPT[23] tree algorithm to cells in

182  reduced dimension space to calculate cell pseudotimes $(\tau_i, ..., \tau_n)$.

183

184  **Selection of TRAINing genes**

185  Because trajectory pseudotime $\tau$ is derived from underlying gene expression

186  profiles, we hypothesized that a trajectory can sufficiently be described by several

187  trajectory informative TRAINing genes: driver genes whose expression levels exhibit

188  strong linear relationships with pseudotime, and presumably have a greater influence

189  on pseudotime calculation and graph learning. Biologically, we assume that genes

190  with strong linear relationships with pseudotime are highly significant in

191  differentiation processes. Specifically, consider a single trajectory branch **B,**

192  consisting of an $n$ cells by $|G|$ genes expression matrix:

6

193

$$\boldsymbol{B} = \begin{bmatrix} x_{1,A} & \cdots & x_{1,|G|} \\ \vdots & \ddots & \vdots \\ x_{n,A} & \cdots & x_{n,|G|} \end{bmatrix},$$

194    where $n$ is the number of cells in $\boldsymbol{B}$, $A$ denotes a gene, and $x_{1,A}$ denotes the

195    expression of gene $A$ in cell 1. Each cell $(1, \ldots, n)$ has a corresponding pseudotime

196    $(\tau_i, \ldots, \tau_n,)$. We aim to identify influential TRAINing genes by using gene-pseudotime

197    covariance as a metric for evaluating gene significance in a differentiation trajectory:

$$\boldsymbol{C} = Cov\big(\boldsymbol{B}, (\tau_i, \ldots, \tau_n,)\big)$$

198    In each branch, genes are ranked by covariance and the lowest ranked genes

199    (default: bottom 5%) are removed from the workflow to prevent these from

200    confounding further analysis. The remaining genes are classified as TRAINing genes

201    for that trajectory.

202    We note that TRAINing genes are distinct from commonly used 'differentially

203    expressed genes' in two ways: 1. TRAINing genes are not dependent on cell type

204    annotations, and 2. TRAINing genes may not necessarily exhibit large absolute

205    changes in expression as one traverses a cell lineage but strongly co-vary with

206    pseudotime. It is this covariance, rather than absolute expression, that is used to

207    define TRAINing genes.

208

209    While covariance is the default metric, ENTRAIN can alternatively be configured to

210    use correlation coefficients.

211

212    **Extracting regulatory information from NicheNet**

213    Expression dynamics during differentiation are likely to be a manifestation of cell-

214    intrinsic and cell-extrinsic regulatory programmes. To demarcate these two factors,

215    the algorithm's second step unites prior knowledge of ligand-receptor pairs and their

216    corresponding intracellular regulatory interactions to determine potential ligands

217    driving the observed TRAINing gene expression dynamics.

218

219    Under the assumption that the microenvironmental niche has a quantifiable

220    contribution to gene expression dynamics in differentiation, we require a database

221    that predicts which target genes are subject to regulation by ligand-receptor pairs.

222    ENTRAIN extracts this information from NicheNet [24], which unites traditional ligand-

223    receptor signalling to downstream transcriptional regulation. We first identified active

224 LR pairs amongst the trajectory cells ('receivers') and the remaining cells in the

225 dataset ('senders'), using NicheNet as prior knowledge of possible ligand-receptor

226 interactions. With the assumption that high LR expression levels do not necessarily

227 correlate to significance in driving differentiation trajectories, we determined LR pairs

228 for further analysis if they fulfilled two criteria: 1. They are expressed by a sufficient

229 proportion of cells in the dataset (default >0 counts in at least 10% of cells). 2. The

230 corresponding receptors are expressed by a sufficient proportion of differentiating

231 cells (default >0 counts in at least 10% of differentiating cells). Of the ligands that

232 meet the criteria, we extracted their respective downstream target regulation scores

233 from the NicheNet database. These are vectors representing the ability of a given

234 ligand to regulate every human gene. Thus, each ligand is associated with a vector

235 of length $g$, where $g$ is the number of human genes in the database, and each

236 element of the vector is a number (a "regulatory potential") representing the strength

237 of the regulatory relationship between the ligand and a given gene.

238

239 **Calculation of top environmental drivers of a trajectory**

240 Next, we assumed that some subset of the active ligands will constitute the

241 extracellular signals influencing a trajectory. We speculated that the regulation

242 between ligands and the trajectory could be contained in existing databases of

243 regulatory networks interactions.

244

245 To detect this, we used a supervised random forest model to fit NicheNet regulatory

246 potentials (explanatory variable) to TRAINing gene covariances (response variable)

247 [25]. Here, we consider the NicheNet matrix as an $L$ by $|G|$ matrix $\mathbf{L}$, where $L$ is the

248 number of actively signalling ligands, and the covariances are represented by a $|G|$

249 dimensional column vector $\mathbf{C}$. Random forest attempts to fit $\mathbf{L}$ to $\mathbf{C}$, used with

250 hyperparameters *n_trees* = 500, *n_features at each split* = number of ligands

251 (features) divided by 3.

252

253 In principle, some columns of $\mathbf{L}$ (which represent the predicted change in gene

254 expression as a result of the ligand-receptor pairing), will possess greater similarity

255 to $\mathbf{C}$ than others if the ligand is responsible for the observed covariance in $\mathbf{C}$. This

256 similarity is represented as variable importance, calculated by removing one column

257 at a time from the matrix and calculating the loss in Gini index that results from the

258 removal. Thus, variable importance represents the significance of a ligand in

259 predicting observed gene expression covariance.

260

261 To assess the environmental dependence of whole trajectory branches, we used %

262 Variance Explained (%V.E.). This metric measures how well the random forest

263 predicts the variance in $C$. More formally,

$$\%\text{V.E.} = 1 - \frac{MSE}{Var(\boldsymbol{C})}$$

264

265 Random forest was chosen as the primary algorithm for feature scoring owing to

266 several advantages suited for our context. Firstly, it caters to non-linear interactions

267 between features, such as those that might be found in regulatory interactions

268 between ligands and their downstream target genes. Secondly, built-in methods for

269 feature selection and scoring, based on sequential removal of features,

270 accommodates our primary goal of scoring ligands rather than predicting gene

271 expression. Thirdly, while a known drawback of random forests is the difficulty of

272 interpretability, this is offset by our existing prior knowledge of gene regulatory

273 networks that provides the insight into downstream targets. Lastly, considering the

274 relatively low numbers of ligands and receptor genes relative to the rest of the

275 genome, the computational complexity of random forests compared to other feature

276 selection algorithms becomes less concerning. Moreover, our fitting is performed on

277 the level of trajectory branches or velocity clusters, rather than individual cells,

278 further mitigating concerns of computational complexity.

279

280 **Calculation of cell-wise influences**

281 Differentiating cells exhibit changes in receptor expression and regulatory wiring as

282 they progress along a developmental process. Because of this, we hypothesized that

283 certain stages of a developmental process will be more influenced by environmental

284 signalling than other stages. We thus wished to produce a more granular, cell-wise

285 measure of ligand influence that encapsulates this behaviour. To do this we

286 calculated pseudotime-expression covariances along a rolling window of cells along

287 pseudotime, restricted to separate branches (**Supplementary Algorithm 1**). We

288 used a default window size $w$ and step size $s$ of 10% and 2% of the cells in the

9

289     trajectory branch, respectively. This 'local covariance' quantifies a gene's expression

290     dynamics within a rolling window of differentiating cells. To this end, we fit a second

291     round of random forest models to each rolling window, such that every branch is now

292     subject to an additional 50 'local' model fits corresponding to 50 rolling windows

293     along the branch. The number of local model fits is dependent on the values of $s$ and

294     $w$; 50 rolling windows is the behaviour when $s$ and $w$ are assigned default values.

295     We used regulatory potentials from the top 5 ligands as the predictor variable (a $|G|$ $x$

296     5 matrix with default parameters) and the local covariances as the response variable.

297     For step sizes greater than 1, we linearly interpolate %V.E.$_{.i}$ values for cells which

298     are skipped.

299     Resultant %V.E. values denote the confidence of the NicheNet fit at each of the 50

300     windows. Genes possessing high covariance with pseudotime are assumed to be

301     important for trajectory determination, and we are interested in the subset of those

302     that are under environmental control. Some of these high-covariance genes will not

303     be under extracellular control and consequently exhibit a low %V.E. value when

304     fitted to NicheNet. On the other hand, high covariance genes that are also under

305     extracellular control will exhibit both high covariances and a confident fit (increased

306     %V.E. ) to NicheNet. As a result, these window %V.E. values can be interpreted as

307     the degree of environmental dependence across different stages of the trajectory.

308     Ultimately, every trajectory branch is subject to one 'branch-wide' model fit that

309     determines the top few ligands of interest, and 50 'local' model fits that assess where

310     their regulatory effects are most noticeable. Cells with cell-intrinsic drivers would be

311     expected to exhibit low, negative, or widely varying %V.E. values as the model

312     cannot accurately fit environmental regulators to the observed expression dynamics

313     in that window, while the opposite is true for highly environmentally dependent

314     windows. We note that the term 'cell-wise' is slightly misleading, as the observed

315     expression dynamics are deduced from the covariances of many neighbouring cells

316     in a rolling window of observations rather than a single cell.

317

318     **Finding ligands responsible for RNA velocity dynamics**

319     RNA velocity is a dynamical approach that calculates the time-derivative of RNA

320     concentration for single cells, allowing for short-term predictions of cell fate in

321    differentiating populations. Because these dynamics are often dependent on

322    environmental signals, we predicted that ENTRAIN could be employed to determine

323    driver ligands responsible for observed RNA velocity vectors. Biologically, these

324    represent ligands that may be responsible for short time scale dynamics that may not

325    be resolvable using the pseudotime-based approach described previously.

326    For full details of the velocity estimation, see ref. [26].

327

328    In most datasets, a small minority of genes are responsible for the majority of

329    observed velocity variance[27], necessitating a way to prioritize velocity genes by their

330    significance. The ENTRAIN-Velocity module uses scVelo to recover fit likelihoods, a

331    measure of velocity significance [26], from which to infer ligand activity

332    (**Supplementary Figure S1**).

333

334    We first clustered the RNA velocity matrix into $c$ groups representing major axes of

335    variance in RNA velocity vectors, by repurposing the Leiden algorithm in scanpy[28].

336    We then calculated the fit likelihoods for velocity genes, by applying the scVelo

337    recover_dynamics[26] function to each velocity cluster. For each velocity cluster $c_i$, this

338    process generates a vector $\boldsymbol{\ell}_i$ of length $|G_i|$, where $|G_i|$ is the number of genes with

339    calculated fit likelihoods per cluster $c_i$. Note that the genes with calculated fit

340    likelihoods are usually a subset of all genes because not all genes possess confident

341    velocities. These genes (row names) constitute our TRAINing genes for this module,

342    and the fit likelihoods (values) represent the response variable for subsequent model

343    fit described below.

344

345    To elucidate environmental influence driving the velocities, we fit the NicheNet

346    ligand-target matrix to all genes with calculated likelihoods using a random forest

347    regression model [25] with hyperparameters *n_trees* = 500, *n_features at each split* =

348    number of ligands (features) divided by 3. As before, we consider the NicheNet

349    matrix as an $L$ by $|G_c|$ matrix $\boldsymbol{L}$, and the velocity likelihoods for a given cluster $c_i$ are

350    represented by a $|G_c|$ dimensional column vector $\boldsymbol{\ell}_c$. Random forest attempts to fit $\boldsymbol{L}$

351    to $\boldsymbol{\ell}_i$ for all clusters $c$ (**Supplementary Algorithm 2**), under the assumption that if a

352    ligand is truly responsible for some component of the observed velocities in a cluster,

353    the corresponding column in $\boldsymbol{L}$ will be more similar to the velocity likelihood vector

11

354 compared to less significant ligands. Similarly to the pseudotime-based approach,
355 we extracted mean decrease in Gini index and %V.E. scores to evaluate ligand
356 significance.

357

358 **Finding ligands responsible for RNA velocity dynamics in spatially resolved**
359 **datasets.**
360 The third module of ENTRAIN, called ENTRAIN-Spatial, is designed for datasets
361 with paired scRNA-seq and Visium data. This module first calculates and clusters
362 velocities on the scRNA-seq matrix object, as in ENTRAIN-Velocity. This is followed
363 by transferring velocity cluster labels to the Visium dataset using the package
364 tangram-sc[29]. Next, within each velocity cluster, the ENTRAIN-Spatial subsets the
365 Visium dataset to include only those spots matching the velocity cluster label or the
366 spots in direct adjacency.

367

368 Subsequently, we select genes that are included in NicheNet's ligand-receptor
369 network to inform later analysis of ligand-receptor pairings. In contrast to the
370 previous ENTRAIN-Velocity module, these genes are restricted to those that are
371 situated in the immediate spatial vicinity of differentiating cells.
372 Subsequent ligand-receptor pairing, random forest fitting, and scoring were
373 performed identically as in the ENTRAIN-Velocity module.

374

375 **RESULTS**
376 ENTRAIN explicitly incorporates output from established trajectory tools to inform a
377 random forest feature selection model for ligand scoring (**Fig. 1A).** As a proof-of-
378 concept, we validated ENTRAIN on a scRNA-seq dataset profiling the bone marrow
379 microenvironment (BME) and its resident mesenchymal and haematopoietic lineages
380 in mice (**Fig. 1B**). We evaluated the contribution of each gene towards the trajectory
381 dynamics by calculating pseudotime using Monocle3 (**Fig. 1C**). We extracted the
382 gene expression for cells along the pre-B trajectory (**Fig. 1D**) and derived the
383 pseudotime-expression covariance for every gene. We assessed the biological
384 relevance of this metric by ranking the genes by covariance and interrogating the top
385 covarying TRAINing genes. This revealed known lineage marker genes including
386 *Vpreb1, Igll1,* and *Vpreb3* for pre-B cells. In parallel, we examined the
387 microenvironmental interactions by selecting receptor and ligand genes. We then

12

388 queried the NicheNet ligand-target regulatory potential database to obtain regulatory

389 interactions between active ligands and their corresponding regulons. ENTRAIN was

390 then performed on the developing B cell lineages using this database as input. We

391 calculated the model's V.E., a measure of the proportion of TRAINing gene

392 covariance that can be attributed to extracellular signals. The percentage of V.E. by

393 the 71 identified active ligands was 2.6%. To identify more granular behaviour, we

394 conducted ENTRAIN in a cell-wise manner by analysing environmental dependence

395 in a series of 100 rolling windows along trajectory pseudotime for every branch. This

396 analysis revealed that the previous environmental dependence was restricted to

397 small pockets of HSCs (**Fig. 1E**), indicating that the local ligand influence was

398 restricted to a subpopulation of progenitor cells that appeared relatively early in

399 lineage commitment. ENTRAIN output shortlisted signalling ligands that are known

400 to be involved in B cell development (*Vcam1/Lama2-Itgb, Il7-Il7r*; *Tnfsf13b-*

401 *Tnfrsf13b; Il15/Il2-Il2rg*) and ligands with conserved roles during cellular

402 differentiation (*Dll1-Notch1/2; Dkk2-Lrp6; Jag1-Notch1/2*) (**Fig. 1F**). The regulatory

403 potential was dominated by a small subset of functionally relevant target genes,

404 particularly *Ebf1, Myl4* and *Cd79a*. Interrogating the source of these ligands revealed

405 that while some of the top-ranked ligands were expressed primarily by a singular cell

406 type (**Fig. 1F**, coloured lines), others were expressed among heterogenous cell

407 types (grey lines). ENTRAIN also identified a novel extracellular signal that was not

408 previously known to be involved in B cell development (*Ptdss1-Scarb1/Jmjd6*) (**Fig.**

409 **1F**).

410

411 To demonstrate the versatility of ENTRAIN we developed the ENTRAIN-Velocity

412 module to recover environmental signals responsible for the RNA velocity vector and

413 applied it to a murine embryonic neurogenesis dataset[30] (**Figure 2A)**. The velocity

414 matrix was recovered using scVelo and clustered with the Leiden algorithm[31] to

415 deconvolute velocity variance into major groups. The vectors formed 10 velocity

416 clusters (VC0-9), which roughly corresponded to major cell types and transitions

417 (**Figure 2B**). We analysed and ranked the joint likelihoods of the velocities in each

418 cluster to identify the TRAINing genes for this dataset: the most rapidly up- or down-

419 regulated genes during neurogenesis (**Fig. 2C**). We then applied ENTRAIN to each

420 velocity cluster to identify driver ligands responsible for the observed velocities (**Fig.**

421 **2D**). The analysis predicted positive V.E. scores for 5 out of 10 clusters (VC0-VC3

422    and VC7) corresponding to velocities exhibited by fibroblastic, radial glial,

423    neuroblast/neuronal, and neural tube cell clusters (**Fig. 2E**). In these clusters, the

424    environmental influence was attributed to ligands in the *Notch* pathway (*Tgfb2,*

425    *Bmp2, Ntf3* and *Bdnf*) and *Wnt* signaling pathway (*Sema3b, Psap* and *Pdgfb*) known

426    to be involved in embryonic neurogenesis. More generally, we considered ligands

427    ranked among the top 5 in each positive cluster and showed that 21 out of 25 ligands

428    were known to be involved in embryonic neurogenesis (**Supplementary Table S1**),

429    with the exceptions being the extracellular matrix proteins *Npnt*/*Adam15* and

430    *Serpinc1.* Interrogation of the NicheNet ligand-target network revealed interactions

431    between *Tgfb2-Ina*/*Mapt*/S*tmn2/Igfbpl1, Bdnf-Bcl11b, and Jag1-Ebf1* as major

432    components of environment-driven neuronal differentiation (vcluster1 and vcluster7),

433    as well as *Jag1-Sdc2* as the largest environmental driver in mesenchymal

434    development (vcluster2) (**Fig. 2F**). Fibroblasts and neuroblasts were the major cell

435    types responsible for producing the highest 3 ranked ligands (**Fig. 2F**).

436

437    Emerging spatial transcriptomics technologies have recently shown success in

438    delineating the role of cell-cell communication in various cellular contexts[21, 32].

439    Building upon this, we developed the ENTRAIN-Spatial module to decode cell-cell

440    communication signals driving RNA velocities, while concurrently considering their

441    spatial environment. This module operates by accepting a paired dataset of spatial

442    transcriptomics data and single-cell data. Its output comprises those ligand-receptor

443    pairs that are both spatially co-localized and have a quantifiable influence on the

444    observed RNA velocities.

445

446    We applied ENTRAIN-Spatial to a paired dataset consisting of both 10x Chromium

447    single-cell and 10x Visium data, which was obtained from Ratz et al.[33] (**Fig. 3A** and

448    **Fig. 3B**). We recovered the RNA velocities from the 10x Chromium data using

449    scVelo[26] and subsequently clustered the velocities into 8 major clusters (**Fig. 3C**). By

450    utilizing Tangram[29], we transferred the velocity cluster labels to their spatial

451    positions.

452    We then used ENTRAIN-Spatial to evaluate ligands located in close spatial proximity

453    to spots associated with a specific velocity label. The scoring was performed based

454    on each ligand's potential to instigate the observed RNA velocities. ENTRAIN-Spatial

455    results indicated that five out of the eight major velocity clusters (vcluster0, 1, 3, 4

456   and 6) exhibited a detectable level of environmental influence, as quantified by the
457   percentage of variance explained (% V.E.) (**Fig. 3D**). Notably, the velocity cluster
458   corresponding to immature and mature oligodendrocytes (vcluster3) demonstrated
459   the highest proportion of variance explained. This cluster corroborated ligands that
460   are well-documented to be implicated in oligodendrocyte maturation, including the
461   *Wnt*-family and *Vgf*. Notably, as opposed to ENTRAIN-Velocity, these ligands are
462   restricted to those expressed in any spot adjacent to a spot associated with a
463   velocity cluster.
464   To interpret spatial patterns in driver ligand expression, ENTRAIN-Spatial facilitates
465   the visualization of specific spots expressing the highest-ranking ligands (**Fig. 3E**) as
466   well as the relative contributions of spatially adjacent cell types towards driving the
467   observed velocities (**Fig. 3F**).
468
469   To corroborate our findings, we benchmarked the performance of ENTRAIN to
470   similar methods NicheNet[18] and CellCall[19] for single-cell RNA results, and Giotto[34]
471   for spatial transcriptomics results, concentrating specifically on the top 10 ligands
472   from each method, as well as the highest velocity confidence clusters
473   (**Supplementary Figure S2**), to maintain. Despite the observed discrepancy
474   between all these methodologies (**Supplementary Fig. S3**), ENTRAIN
475   demonstrated the highest rate of literature support across the top ranked ligands
476   when analyzing the pre-B cell, neuroblast, and oligodendrocyte lineages (**Fig. 3G**).
477
478   These results indicate that ENTRAIN accurately recovers extracellular regulators
479   that are not resolved by DEG-based methods.
480
481   **DISCUSSION**
482   ENTRAIN uses an orthogonal approach that has several advantages over other
483   methods that are highly dependent on the accurate identification of DEGs, which in
484   turn depend on correct and reproducible cell type clusters, labels and pair-wise
485   comparisons. As a result, these methods cannot consider intra-cluster expression
486   dynamics that may arise as a cell differentiates along a trajectory. In comparison,
487   ENTRAIN can be executed on any arbitrary number of cell states linked by a
488   trajectory or RNA velocity vectors. In turn, ENTRAIN can analyse sparse populations
489   that are not amenable to DEG-based methods.

490

491    ENTRAIN exhibits several limitations. Firstly, ENTRAIN-Pseudotime is dependent on

492    the quality of the topology that is learnt by the trajectory inference algorithm[35]. To

493    mitigate this, the ENTRAIN-Pseudotime module allows flexible input from any

494    trajectory method provided that each input cell is assigned a pseudotime value and a

495    trajectory branch in the Seurat object metadata. In addition, ENTRAIN allows

496    interactive selection of trajectory nodes for flexible analysis on a user-defined

497    branch. Secondly, ENTRAIN-Velocity is similarly subject to the same limitations as

498    RNA velocity. Namely, the potential for inferring spurious velocity vectors when it is

499    applied to populations with multiple kinetic regimes or datasets containing mature

500    cell types missing intermediate cell states[12]. Thirdly, the NicheNet database does not

501    discriminate between up- or down-regulated targets, which may result in ENTRAIN

502    detecting both inhibitors and activators of a differentiation pathway. Lastly, ENTRAIN

503    requires whole-transcriptome based technologies to ensure accurate capture of all

504    ligand and receptor genes. Therefore, hybridization-based technologies which detect

505    a limited panel of  genes may not be suitable.

506

507    In conclusion, we present ENTRAIN, the first tool to date that integrates trajectory

508    and cell-cell communication methods to identify driving ligands influencing cell

509    differentiation. Validating ENTRAIN on existing single-cell pre-B Cell, neuronal, and

510    spatially resolved brain datasets demonstrates that ENTRAIN recovers cell-extrinsic

511    determinants of differentiation. Comparative analysis suggests that ENTRAIN

512    outperforms other cell-cell communication methods in deciphering intercellular

513    signals governing differentiation, possibly owing to the leveraging of trajectory and

514    velocity data rather than traditional differential expression. Future work may consist

515    of extension towards capturing epigenetic contributions from methylation or

516    chromatin accessibility data[36, 37]

517

527

528

529  **References**

530  1.      Armingol, E., Officer, A., Harismendy, O. & Lewis, N.E. Deciphering cell-cell

531  interactions and communication from gene expression. *Nat Rev Genet* **22**, 71-88

532  (2021).

533  2.      Aldridge, S. & Teichmann, S.A. Single cell transcriptomics comes of age.

534  *Nature communications* **11**, 4307 (2020).

535  3.      Nguyen, A., Khoo, W.H., Moran, I., Croucher, P.I. & Phan, T.G. Single Cell

536  RNA Sequencing of Rare Immune Cell Populations. *Frontiers in immunology* **9**, 1553

537  (2018).

538  4.      Moris, N., Pina, C. & Arias, A.M. Transition states and cell fate decisions in

539  epigenetic landscapes. *Nat Rev Genet* **17**, 693-703 (2016).

540  5.      Trapnell, C. Defining cell types and states with single-cell genomics. *Genome

541  research* **25**, 1491-1498 (2015).

542  6.      Saelens, W., Cannoodt, R., Todorov, H. & Saeys, Y. A comparison of single-

543  cell trajectory inference methods. *Nat Biotechnol* **37**, 547-554 (2019).

544  7.      Trapnell, C. et al. The dynamics and regulators of cell fate decisions are

545  revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386

546  (2014).

547  8.      Tritschler, S. et al. Concepts and limitations for learning developmental

548  trajectories from single cell genomics. *Development* **146** (2019).

549  9.      La Manno, G. et al. RNA velocity of single cells. *Nature* **560**, 494-498 (2018).

550  10.     Bergen, V., Lange, M., Peidli, S., Wolf, F.A. & Theis, F.J. Generalizing RNA

551  velocity to transient cell states through dynamical modeling. *Nat Biotechnol* **38**,

552  1408-1414 (2020).

553  11.     Bergen, V., Soldatov, R.A., Kharchenko, P.V. & Theis, F.J. RNA velocity-

554  current challenges and future perspectives. *Mol Syst Biol* **17**, e10282 (2021).

555  12.     Alquicira-Hernandez, J., Powell, J.E. & Phan, T.G. No evidence that

556  plasmablasts transdifferentiate into developing neutrophils in severe COVID-19

557  disease. *Clin Transl Immunology* **10**, e1308 (2021).

558    13.    Weinreb, C., Wolock, S., Tusi, B.K., Socolovsky, M. & Klein, A.M.
559    Fundamental limits on dynamic inference from single-cell snapshots. *Proceedings of*
560    *the National Academy of Sciences of the United States of America* **115**, E2467-
561    E2476 (2018).
562    14.    Pavlicev, M. et al. Single-cell transcriptomics of the human placenta: inferring
563    the cell communication network of the maternal-fetal interface. *Genome research* **27**,
564    349-361 (2017).
565    15.    Almet, A.A., Cang, Z., Jin, S. & Nie, Q. The landscape of cell-cell
566    communication through single-cell transcriptomics. *Curr Opin Syst Biol* **26**, 12-23
567    (2021).
568    16.    Efremova, M., Vento-Tormo, M., Teichmann, S.A. & Vento-Tormo, R.
569    CellPhoneDB: inferring cell–cell communication from combined expression of multi-
570    subunit ligand–receptor complexes. *Nature Protocols* **15**, 1484-1506 (2020).
571    17.    Vento-Tormo, R. et al. Single-cell reconstruction of the early maternal–fetal
572    interface in humans. *Nature* **563**, 347-353 (2018).
573    18.    Browaeys, R., Saelens, W. & Saeys, Y. NicheNet: modeling intercellular
574    communication by linking ligands to target genes. *Nature methods* **17**, 159-162
575    (2020).
576    19.    Zhang, Y. et al. CellCall: integrating paired ligand-receptor and transcription
577    factor activities for cell-cell communication. *Nucleic Acids Res* **49**, 8520-8534 (2021).
578    20.    Joost, S. et al. Single-Cell Transcriptomics Reveals that Differentiation and
579    Spatial Signatures Shape Epidermal and Hair Follicle Heterogeneity. *Cell Syst* **3**,
580    221-237.e229 (2016).
581    21.    Almet, A.A., Cang, Z., Jin, S. & Nie, Q. The landscape of cell–cell
582    communication through single-cell transcriptomics. *Current Opinion in Systems*
583    *Biology* **26**, 12-23 (2021).
584    22.    Cao, J. et al. The single-cell transcriptional landscape of mammalian
585    organogenesis. *Nature* **566**, 496-502 (2019).
586    23.    Mao, Q., Yang, L., Wang, L., Goodison, S. & Sun, Y.  792-800 (2015).
587    24.    Browaeys, R., Saelens, W. & Saeys, Y. NicheNet: modeling intercellular
588    communication by linking ligands to target genes. *Nature Methods* **17**, 159-162
589    (2020).
590    25.    Breiman, L. Random Forests. *Machine Learning* **45**, 5-32 (2001).

591    26.    Bergen, V., Lange, M., Peidli, S., Wolf, F.A. & Theis, F.J. Generalizing RNA

592    velocity to transient cell states through dynamical modeling. *Nature Biotechnology*

593    **38**, 1408-1414 (2020).

594    27.    Bergen, V., Soldatov, R.A., Kharchenko, P.V. & Theis, F.J. RNA velocity-

595    current challenges and future perspectives. *Molecular systems biology* **17**, e10282-

596    e10282 (2021).

597    28.    Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene

598    expression data analysis. *Genome Biology* **19**, 15 (2018).

599    29.    Biancalani, T. et al. Deep learning and alignment of spatially resolved single-

600    cell transcriptomes with Tangram. *Nature Methods* **18**, 1352-1362 (2021).

601    30.    La Manno, G. et al. Molecular architecture of the developing mouse brain.

602    *Nature* **596**, 92-96 (2021).

603    31.    Traag, V.A., Waltman, L. & van Eck, N.J. From Louvain to Leiden:

604    guaranteeing well-connected communities. *Scientific reports* **9**, 5233 (2019).

605    32.    Longo, S.K., Guo, M.G., Ji, A.L. & Khavari, P.A. Integrating single-cell and

606    spatial transcriptomics to elucidate intercellular tissue dynamics. *Nature Reviews*

607    *Genetics* **22**, 627-644 (2021).

608    33.    Ratz, M. et al. Clonal relations in the mouse brain revealed by single-cell and

609    spatial transcriptomics. *Nature Neuroscience* **25**, 285-294 (2022).

610    34.    Dries, R. et al. Giotto: a toolbox for integrative analysis and visualization of

611    spatial expression data. *Genome Biology* **22**, 78 (2021).

612    35.    Saelens, W., Cannoodt, R., Todorov, H. & Saeys, Y. A comparison of single-

613    cell trajectory inference methods. *Nature Biotechnology* **37**, 547-554 (2019).

614    36.    Tedesco, M. et al. Chromatin Velocity reveals epigenetic dynamics by single-

615    cell profiling of heterochromatin and euchromatin. *Nature Biotechnology* **40**, 235-244

616    (2022).

617    37.    Vandereyken, K., Sifrim, A., Thienpont, B. & Voet, T. Methods and

618    applications for single-cell and spatial multi-omics. *Nature Reviews Genetics* (2023).
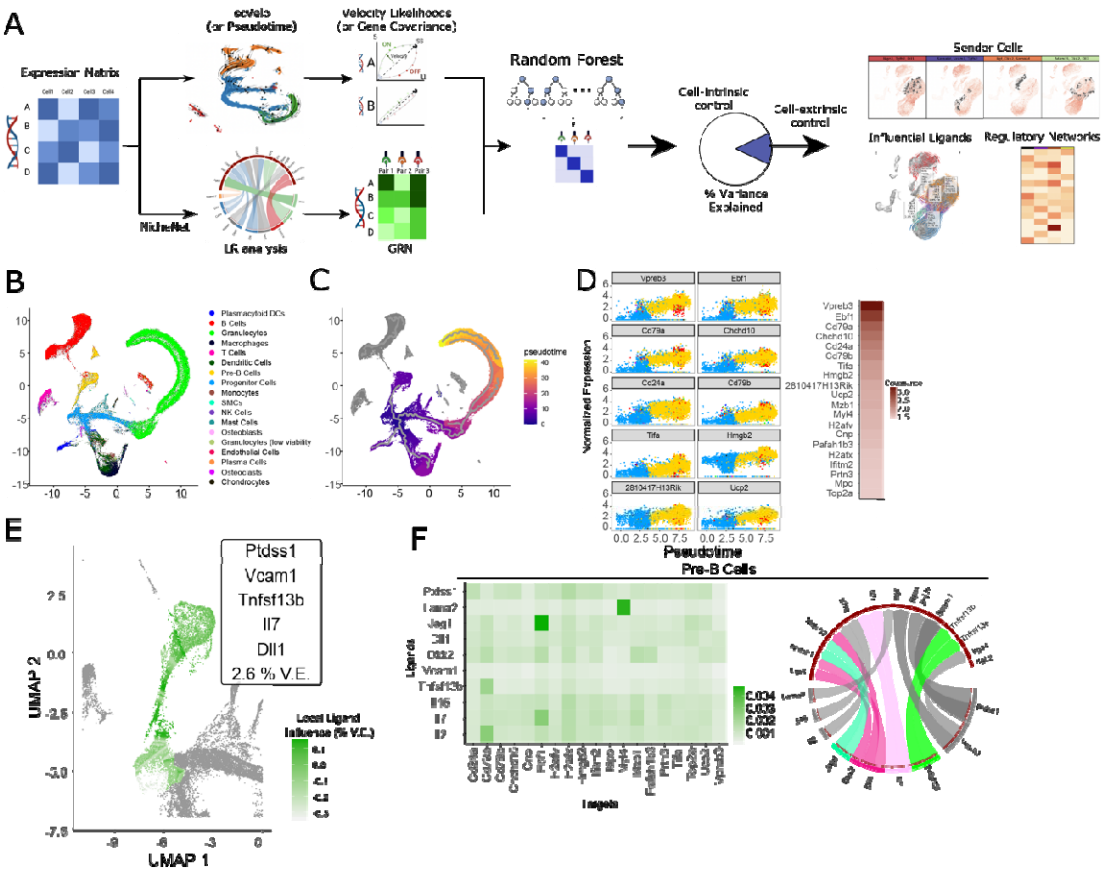
619

620

621

622

**FIGURE 1: ENTRAIN-Pseudotime analysis of pre-B cell development.**

(A) ENTRAIN workflow.

(B) UMAP representation of 133,942 cells in mouse bone marrow environment.

(C) Monocle3 trajectory overlayed on the UMAP.

(D) High trajectory covariance (TRAINing) genes for the trajectory between haematopoietic progenitors and pre-B cells.

(E) ENTRAIN ligand results overlayed on the B cell lineage trajectory. Cells coloured by local ligand influence. V.E: Variance Explained

(F) Ligand-target gene regulatory networks (left) and circos plot (right) representing regulatory links between top ranked ligands and their downstream targets. Colour represents identity of major cell type expressing that ligand. Ligands expressed by more than one cell type are coloured grey.
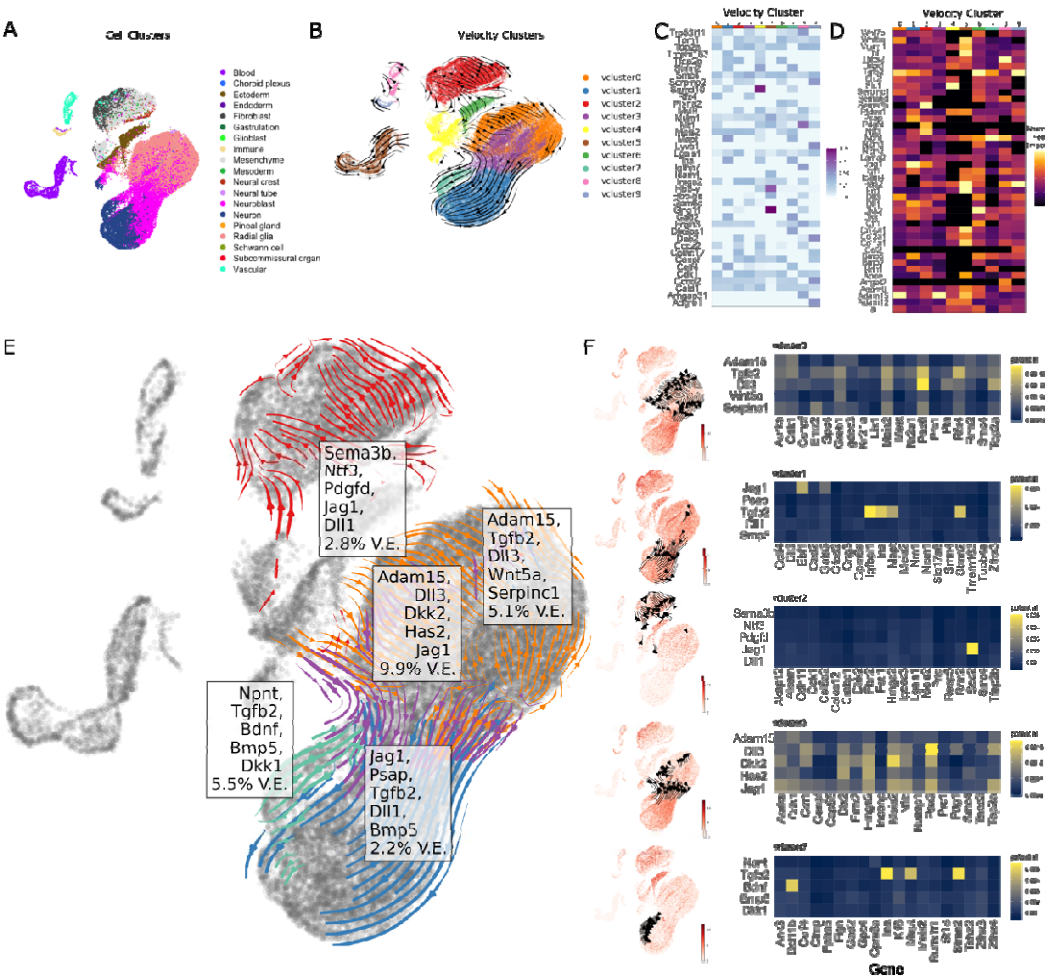
**FIGURE 2: ENTRAIN-Velocity analysis of neuronal development.**

(A) UMAP representation of mouse embryonic neurogenesis dataset at E10.5.

(B) Velocity vectors overlayed on UMAP representation, cells coloured by velocity cluster membership.

(C) Heatmap of high likelihood velocity genes in each velocity cluster.

(D) Heatmap of ligands predicted by ENTRAIN to influence velocities in each velocity cluster.

(E) Velocity vectors, V.E. scores and top 5 ligands predicted by ENTRAIN for 5 out of 10 clusters (VC0-3, VC7) overlayed on the UMAP embedding.

(F) Sender expression and predicted gene targets for the top 5 ligands in each velocity cluster. Left: UMAP embedding coloured by mean expression of the top 5 ligands predicted for the velocity cluster. Right: Heatmap showing NicheNet regulatory linkages between the top 5 ligands (y-axis) and their downstream target genes (x-axis).e
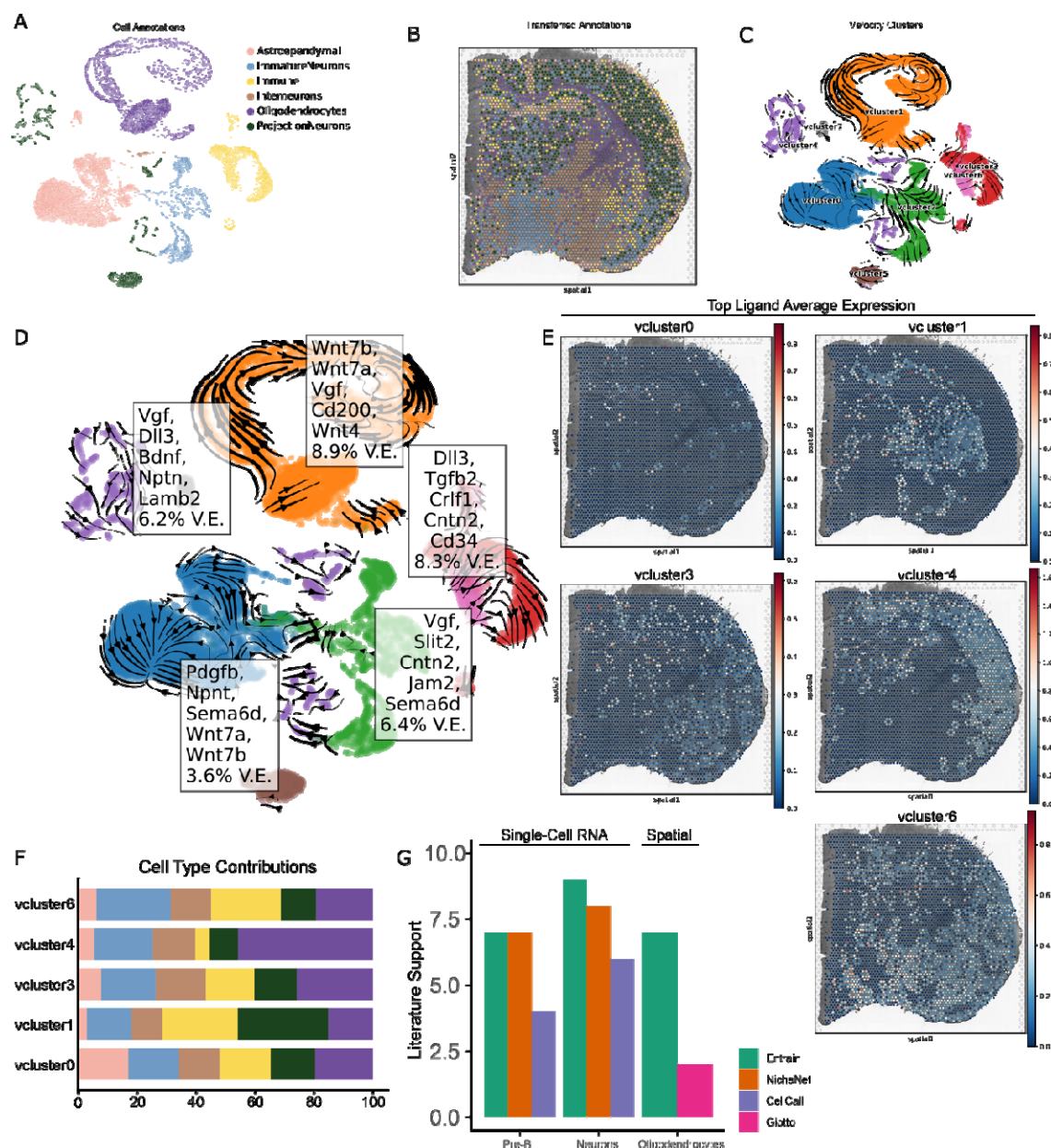
**FIGURE 3: ENTRAIN-Spatial analysis of neuronal development at spatial resolution.**

(A) UMAP plot of pre-annotated Ratz et al. dataset

(B) Tangram transferred labels overlayed on spatial scatter plot.

(C) UMAP plot of velocity cluster labels.

(D) Top 5 ligands predicted by ENTRAIN for positive V.E. clusters (Velocity Clusters 0, 1, 3, 4 and 6) overlayed on velocity plot.

(E) Spatial scatter plot representing average expression of top 5 ligands associated with each velocity cluster.

660    (F) Stacked bar plot showing the proportion of cell types expressing the top 5 ligands

661         for each velocity cluster, weighted by variance explained.

662    (G) Bar plot showing number of ligands with literature support for their role in pre-B

663         cell and neuronal development.